ABSTRACT OF THE INVENTION

The present invention relates to methods for identifying novel genes comprising: (i) generating one or more specialized databases containing information on gene/protein structure, function and/or regulatory interactions; and (ii) searching the specialized databases for homology or for a particular motif and thereby identifying a putative novel gene of interest. The invention may further comprise performing simulation and hypothesis testing to identify or confirm that the putative gene is a novel gene of interest. The present invention also relates to natural language processing and extraction of relational information associated with genes and proteins that are found in genomics journal articles. To enable access to information in textual form, the natural language processing system of the present invention provides a method for extracting and structuring information found in the literature in a form appropriate for subsequent applications.